

## **Tugas NLP**

1. Judul Artikel : An NLP-Inspired Data Augmentation Method for Adverse Event Prediction Using an Imbalanced Healthcare Dataset
2. Penulis : TOMOKI ISHIKAWA, TAKAHIRO YAKOH AND HISASHI URUSHIHARA
3. Tahun Publikasi : 2022
4. Nama Jurnal : IEEE
5. DOI : Digital Object Identifier 10.1109/ACCESS.2022.3195212

### **Ringkasan Artikel :**

Penelitian ini bertujuan untuk mengatasi masalah ketidakseimbangan data dalam prediksi kejadian merugikan (Adverse Events atau AE) yang timbul dari penggunaan obat. Data medis yang tidak seimbang merupakan tantangan utama dalam analisis prediktif karena dapat mengurangi akurasi prediksi, khususnya untuk kejadian langka yang jarang muncul. Untuk mengatasi hal ini, peneliti mengusulkan pendekatan baru yang terinspirasi dari metode augmentasi data di bidang Natural Language Processing (NLP). Metode ini menghasilkan data sintetis dengan menggantikan informasi latar belakang pasien menggunakan pasangan serupa yang dihitung berdasarkan cosine similarity. Proses dimulai dengan membentuk representasi distribusi dari latar belakang pasien menggunakan model skip-gram, di mana pasangan latar belakang pasien dengan kemiripan tinggi didefinisikan dalam bentuk semacam kamus (thesaurus). Catatan pasien sintetis kemudian dihasilkan dengan mengganti informasi latar belakang asli dengan informasi yang serupa, membentuk dataset yang lebih seimbang. Eksperimen dilakukan pada dataset medis nyata yang terdiri lebih dari 1,5 juta catatan, dan hasilnya menunjukkan bahwa metode augmentasi ini meningkatkan akurasi prediksi hingga 80%, serta F1-score hingga 40% dibandingkan metode yang tidak menggunakan augmentasi. Peningkatan ini terlihat khususnya pada prediksi AE dengan rasio positif rendah (antara 1% hingga 2,1%), yang sering sulit diprediksi dalam skenario medis. Temuan ini menunjukkan efektivitas metode yang diusulkan dalam menangani ketidakseimbangan data pada dataset medis skala besar, sehingga memungkinkan prediksi AE yang lebih akurat dan dapat diandalkan dalam analisis data medis.

### **Analisis Kontribusi Penelitian :**

Penelitian ini memperkenalkan metode augmentasi data yang inovatif untuk prediksi kejadian merugikan di bidang medis, yang terinspirasi dari teknik NLP. Metode yang dikembangkan memiliki pendekatan unik dengan memanfaatkan informasi demografis pasien, yang biasanya diabaikan dalam prediksi kejadian merugikan. Teknik ini menggunakan model skip-gram untuk membentuk pasangan latar belakang pasien serupa berdasarkan nilai kesamaan (cosine similarity), dan selanjutnya menggantikan data asli dengan data sintetis yang mirip. Dalam konteks prediksi AE, inovasi ini memecahkan permasalahan data tidak seimbang yang sering kali mengakibatkan performa prediksi rendah untuk kejadian langka. Metode ini menawarkan solusi untuk meningkatkan performa prediksi tanpa memerlukan data nyata tambahan, yang sering sulit didapatkan atau terbatas dalam konteks medis. Metode ini juga memungkinkan integrasi informasi latar belakang pasien, seperti usia dan riwayat komplikasi, yang jarang diperhitungkan dalam model prediktif tradisional. Temuan ini membawa

kontribusi signifikan pada prediksi AE berbasis machine learning, yang secara tradisional mengabaikan variabilitas latar belakang pasien. Dengan pendekatan augmentasi data berbasis NLP ini, penelitian ini membuka jalan untuk aplikasi yang lebih luas di bidang medis dan memperlihatkan potensi signifikan untuk meningkatkan presisi dan akurasi dalam prediksi berbasis data yang tidak seimbang.

Kelebihan dan Keterbatasan :

Penelitian ini memiliki berbagai kelebihan, terutama dari segi metode augmentasi data yang inovatif dan efektif dalam menangani ketidakseimbangan dataset, masalah umum dalam data medis. Pendekatan yang digunakan sangat efektif dalam menghasilkan data sintetis yang serupa dengan data aktual, sehingga mampu meningkatkan performa prediktif AE secara signifikan. Penelitian ini juga memanfaatkan dataset yang sangat besar, terdiri dari lebih dari 1,5 juta catatan pasien, menunjukkan skalabilitas metode augmentasi ini dalam skenario dunia nyata. Selain itu, model ini memungkinkan pelibatan informasi latar belakang pasien yang jarang dipertimbangkan dalam studi prediksi medis lainnya, seperti usia dan jenis kelamin, sehingga lebih mampu menangkap kompleksitas kejadian medis. Namun, terdapat beberapa keterbatasan dalam penelitian ini. Pertama, fitur data yang digunakan terbatas pada empat jenis informasi latar belakang (usia, jenis kelamin, komplikasi, dan obat yang diberikan), yang mungkin tidak mencakup variabel lain yang dapat meningkatkan akurasi prediksi lebih jauh. Kedua, penelitian ini menggunakan klasifikasi menengah dari kode ICD-10 untuk merepresentasikan penyakit, yang membatasi spesifisitas diagnosis. Semakin mendalam kategori kode ICD-10 yang digunakan, semakin spesifik pula data sintetis yang dihasilkan. Keterbatasan dalam evaluasi juga menjadi tantangan, mengingat metode augmentasi ini dirancang khusus untuk data berbasis teks. Dalam penelitian selanjutnya, mengintegrasikan lebih banyak variabel latar belakang dan metode evaluasi yang lebih komprehensif dapat membantu meningkatkan performa prediksi.

Relevansi untuk NLP Saat Ini :

Hasil penelitian ini dapat diterapkan dalam teknologi NLP yang berurusan dengan data tidak seimbang, khususnya dalam prediksi risiko dalam layanan kesehatan. Sebagai contoh, metode augmentasi ini dapat diterapkan dalam sistem peringatan dini untuk mendeteksi komplikasi pasien di ruang ICU berdasarkan catatan medis elektronik (Electronic Health Records atau EHR). Dengan menerapkan metode ini pada EHR berbasis teks, seperti memo medis atau catatan dokter, prediksi menjadi lebih presisi dengan memanfaatkan data sintetis yang mempertimbangkan informasi latar belakang pasien. Di luar skenario medis, pendekatan ini juga memiliki potensi besar untuk diadaptasi dalam berbagai aplikasi NLP yang menangani data tidak seimbang. Sebagai contoh, dalam analisis teks yang bertujuan untuk klasifikasi otomatis dan pengelompokan risiko pasien, hasil penelitian ini bisa membantu menyeimbangkan dataset dan meningkatkan akurasi model. Selain itu, metode ini dapat dikembangkan lebih jauh untuk aplikasi lain di bidang NLP, seperti klasifikasi diagnosis medis atau analisis sentimen dalam layanan kesehatan. Hasil penelitian ini juga bisa digunakan dalam pengembangan chatbot medis, yang dapat memberikan rekomendasi berbasis risiko kepada pasien. Dengan penyesuaian dan pengembangan lebih lanjut, pendekatan augmentasi data ini dapat memberikan kontribusi signifikan terhadap peningkatan performa model NLP yang berfokus pada prediksi risiko atau klasifikasi pada data tidak seimbang dalam layanan kesehatan.

## Pendapat dan Rekomendasi :

Menurut saya, artikel ini memberikan wawasan baru yang signifikan terkait penerapan teknik NLP untuk mengatasi ketidakseimbangan data dalam prediksi kejadian merugikan pada dataset medis. Penelitian ini tidak hanya memberikan solusi teknis untuk prediksi AE pada dataset yang tidak seimbang, tetapi juga menunjukkan cara inovatif dalam memanfaatkan data latar belakang pasien. Saya percaya bahwa penelitian ini membuka jalan bagi studi lanjutan yang lebih luas, khususnya untuk prediksi kejadian medis dengan menggunakan data sintetis. Namun, untuk memperluas cakupan dan manfaat penelitian ini, penelitian selanjutnya bisa menambahkan lebih banyak variabel latar belakang pasien, seperti kondisi lingkungan atau riwayat keluarga, yang dapat memperkaya model dan meningkatkan akurasi prediksi. Di sisi lain, mengintegrasikan metode augmentasi data ini dengan model berbasis transformer yang lebih kompleks bisa menjadi arah penelitian menarik untuk mengeksplorasi kemampuan prediksi di data medis yang lebih heterogen. Metode yang diusulkan juga bisa diuji pada domain lain di luar kesehatan untuk mengatasi ketidakseimbangan data, misalnya dalam klasifikasi pelanggan untuk industri jasa atau analisis sentimen pada media sosial. Secara keseluruhan, artikel ini memperlihatkan potensi besar dan aplikasi luas dari pendekatan augmentasi data yang diusulkan dalam berbagai bidang, terutama yang berhadapan dengan ketidakseimbangan data.